

1 **DIRECTLY ADDRESSED MULTICAST PROTOCOL**

2 **TECHNICAL FIELD**

3 This invention relates to a new internetworking protocol (IP) multicast protocol
4 that will be useful for implementing networked storage.

5 **BACKGROUND ART**

6 Existing IP multicast technology works by having the sending systems address a
7 packet of data to a unique multicast address, which is then routed by the network
8 infrastructure to each of the remote destinations that have expressed a desire to receive
9 the data packet. As such, the destinations for the IP multicast data packets are unknown
10 to the sender. This model of addressing is useful for the most prevalent use of the current
11 IP multicast scheme, streaming multimedia.

12 In existing IP multicast protocols, the sender sends a data packet to a single virtual
13 multicast address, where the actual recipients are determined by a subscription process
14 managed by a network of switches and routers. There are two alternatives in the existing
15 art for accomplishing multicast data transmission over an IP network: (1) sending
16 multiple singly-addressed packets, one to each of the remote destinations; and (2) setting
17 up a multicast group in the network switch, sending the multicast packet, and then
18 deleting the multicast group. The first alternative results in a significant load on the
19 sending system and on the network infrastructure connected to it. The second alternative
20 results in an increased burden on the network switches. A need exists for an improved
21 method of sending multicast data packets.

22 **SUMMARY OF INVENTION**

23 The present invention uses an enhanced addressing model for IP multicast
24 protocols. In one respect, what is described is a method for transmitting data packets over
25 a network to selected multiple remote destinations wherein the data packets do not travel
26 over any segment of the network more than once and the selected multiple remote
27 destinations are not required to initiate delivery of the data packets through a subscription
28 service. The method includes the following steps: embedding in a data packet a list of
29 multiple remote destination addresses corresponding to each of the selected multiple
30 remote destinations; providing an addressing protocol by which networking elements are
31 capable of accessing the list of multiple remote destination addresses; and instructing the
32 networking elements to transmit a copy of the data packet to each of the selected

1 multiple remote destinations corresponding to each of the addresses in the list of multiple
2 remote destination addresses.

3 In another respect, what is described is a system for transmitting data packets over
4 a network to selected multiple remote destinations wherein the data packets do not travel
5 over any segment of the network more than once and the selected multiple remote
6 destinations are not required to initiate delivery of the data packets through a subscription
7 service. The system includes: a directly addressed multicast protocol (DAMP) client; a
8 network infrastructure; multiple network devices remotely located from the DAMP client
9 and the network infrastructure; and data packets for carrying data between the DAMP
10 client and the multiple network devices.

11 In yet another respect, what is described is a computer-readable medium on which
12 is embedded a program. The embedded program includes instructions for executing the
13 above method.

14 Those skilled in the art will appreciate these and other advantages and benefits of
15 various embodiments of the invention upon reading the following detailed description of
16 an embodiment with reference to the below-listed drawings.

17 BRIEF DESCRIPTION OF DRAWINGS

18 Figure 1 is a block diagram of a system for delivering multicast data packets over
19 a network according to the prior art;

20 Figure 2 is a block diagram of one embodiment of a system for transmitting data
21 packets over a network to selected multiple remote destinations according to the present
22 invention;

23 Figure 3 is a block diagram of another embodiment of a system for transmitting
24 data packets over a network to selected multiple remote destinations according to the
25 present invention;

26 Figure 4 is a flowchart of one embodiment of a method for transmitting data
27 packets over a network to selected multiple remote destinations;

28 Figure 5 is a diagram of one embodiment of the structure of a directly addressed
29 multicast protocol data packet according to the present invention;

30 Figure 6 is a diagram of one embodiment of the structure of an IP options field

1 within a directly addressed multicast protocol data packet according to the present
2 invention;

3 Figure 7 is a diagram illustrating a path of one embodiment of a directly
4 addressed multicast protocol data packet through the Internet to multiple destinations; and

5 Figure 8 is a diagram illustrating another path of one embodiment of a directly
6 addressed multicast protocol data packet through the Internet to multiple destinations.

7 DETAILED DESCRIPTION

8 The present invention relates to a new IP multicast protocol hereinafter referred to
9 as Directly Addressed Multicast Protocol (DAMP). In a DAMP system, a sending client
10 specifies a list of remote destinations directly for each data packet being sent. The term
11 DAMP is used as a label only, and other terms can define the same or similar protocols.

12 Using DAMP, a network of switches and routers will continue to resend the data
13 packets over only those network segments that contain a route to at least one of the
14 specified remote destination addresses. This method thus does not generate unnecessary
15 traffic on any network segment. This form of multicast also greatly simplifies the
16 multicast process by removing the need for multicast group management protocols to
17 manage the destination lists within the network infrastructure. The list of remote
18 destinations are embedded in each data packet, rather than being persistent within the
19 network infrastructure.

20 Figure 1 is a block diagram of a system 100 for delivering multicast data packets
21 over a network according to the prior art. In the existing art covering IP multicast
22 protocols, the system 100 includes at least one IP multicast sending client 110 connected
23 to and communicating with a network infrastructure 115. The IP multicast sending client
24 110 sends IP multicast data packets out to the network infrastructure 115. Within the
25 network infrastructure 115, at least one IP multicast subscription manager 120 is further
26 connected to and communicating with a number of network devices 140. The IP
27 multicast subscription manager 120 receives requests from the network devices 140 to
28 subscribe to multicast data transmissions from the IP multicast sending client 110. When
29 the IP multicast sending client 110 sends IP multicast data packets out to the network
30 infrastructure 115, the IP multicast subscription manager 120 determines which of the
31 network devices 140 have subscribed to those particular data packets. Alongside the IP
32 multicast subscription manager 120 within the network infrastructure 115 are a

1 number of network switches and routers 130. Upon a determination of which network
2 devices 140 are subscribed to the data packets being transmitted from the IP multicast
3 sending client, the network switches and routers 130 forward copies of the data packets to
4 each of the network devices 140 subscribed to the data packets, and await receipt
5 acknowledgments from the network devices 140. The network infrastructure 115 may
6 resend data packets until all data packets are confirmed received by all network devices
7 140. As a result of this multicast scheme, the data packets may travel over any given
8 network segment many times, multiplying the load on the network bandwidth.

9 Figure 2 is a block diagram of one embodiment of a system 200 according to the
10 present invention for transmitting data packets over a network to selected multiple remote
11 destinations. The system 200 of the present invention includes a DAMP sending client
12 210 which is connected to one or more network elements 220. The DAMP sending client
13 210 transmits DAMP data packets to the network elements 220. A network element 220
14 may be a network switch or router and is a part of the overall network infrastructure 115,
15 as described for Figure 1. However, unlike existing multicast systems 100, the present
16 invention does not require an IP multicast subscription manager 120 to manage requests
17 for subscriptions to multicast data transmissions and to manage lists of requesting
18 network devices 140. The system 200 of an embodiment of the present invention
19 includes a number of remote network devices 240. In the present invention, a list of
20 selected remote network devices 240 that are to receive the DAMP multicast transmission
21 is embedded in the DAMP data packets themselves. The network elements 220 route the
22 DAMP data packets on to the selected remote network devices 240 determined in the list
23 embedded in the DAMP data packets. In the system 200, it is not necessary for the
24 DAMP data packets to travel over any given network segment more than once, thus
25 reducing the load on the bandwidth of the network.

26 Figure 3 is a block diagram of another embodiment of a system 300 for
27 transmitting data packets over a network to selected multiple remote destinations
28 according to the present invention. The system 300 illustrates the use of the present
29 invention in a network storage application. The system 300 includes a DAMP sending
30 client 310 connected to one or more network elements 320. The DAMP sending client
31 310 transmits DAMP data packets to the network elements 320. As with the network
32 elements 220 of the system 200, the network elements 320 may include network switches,
33 routers, or other network infrastructure elements. The network elements

1 320 are connected to a number of remote network storage devices 340. As with the
2 embodiment described for Figure 2, in the system 300, a list of selected remote network
3 storage devices 340 that are to receive the DAMP multicast transmission is embedded in
4 the DAMP data packets. The network elements 320 route the DAMP data packets on to
5 the selected remote network storage devices 340 determined in the list embedded in the
6 DAMP data packets.

7 Figure 4 is a flowchart of one embodiment of a method 400 for transmitting data
8 packets over a network to selected multiple remote destinations. The method 400
9 includes the steps of: embedding a list of remote destination addresses into data packets
10 (step 420); enabling network elements to access the list of remote destination addresses
11 (step 430); and instructing network elements to transmit copies of the data packets to each
12 address on the list of remote destination addresses (step 440).

13 In one embodiment of the method 400, the embedding step 420 includes the
14 additional steps of: setting up an IP Options field within the IP header section of a
15 DAMP data packet; setting a Code byte within the IP Options field to a specific value to
16 indicate that the data packet is a DAMP data packet; setting a Length byte to a
17 determinable value to indicate the length in 32-bit words of the IP Options field;
18 embedding in successive 32-bit words the values of a determinable number of IP
19 addresses for the multiple remote destinations for the DAMP data packet; and setting the
20 source IP address in the header section of the DAMP data packet to the IP address of the
21 DAMP sending client 210; and setting the destination IP address in the IP header section
22 of the DAMP data packet to the IP address of one of the multiple remote destinations
23 embedded in the IP Options field. In this embodiment of the method 400, the copying
24 step 440 further includes the steps of: the network element 220 receiving a copy of the
25 DAMP data packet from the DAMP sending client 210 or from another network element
26 220; the network element 220 zeroing those IP addresses embedded in the IP Options
27 field that are not directly accessible below the network element 220; the network element
28 220 setting the destination IP address in the IP header section of the DAMP data packet to
29 the IP address of one of the non-zeroed multiple remote destinations embedded in the IP
30 Options field; and the network element 220 routing a copy of the modified DAMP data
31 packet to each additional network element 220 or network device 240 for which there is a
32 corresponding non-zeroed IP address listed in the embedded list of multiple remote

1 destination IP addresses.

2 Figure 5 is a diagram of the structure 500 of one embodiment of a DAMP data
3 packet 510 according to one embodiment of the present invention. The data packet
4 structure 500 includes a data packet header section 515 and a data packet data section
5 525.

6 The data packet header section 515 contains multiple fields, one of which is a
7 variable length field containing certain IP options information, an IP options field 520.
8 Figure 6 is a diagram of one embodiment of the structure of an IP options field 520 within
9 a directly addressed multicast protocol data packet 510 according to the present invention.
10 The variable length IP options field 520 comprises a sequence of items, each of which
11 start with a Code byte 610 comprised of the following bits, where bit 0 is considered the
12 Most Significant Bit ("MSB"):

13 Bit 0 - Copy bit - which defines whether the IP options field 520 should be copied
14 into each network fragment if the data packet 510 is fragmented, or split, across
15 multiple network frames. A value of 0 indicates that the IP options field 520
16 should be copied into only the first frame, whereas a value of 1 indicates that the
17 IP options field 520 should be copied into every frame;

18 Bits 1-2 - Option class - which defines a set of IP class values;

19 Bits 3-7 - Option number - these bits identify the specific IP option for this data
20 packet, where each of the available IP options is associated with a unique Option
21 number.

22 An IP data packet 510 may be encoded with DAMP using an IP option Code byte
23 610 which in one embodiment may be set to a value of 138. Other values for the IP
24 option Code byte 610 may be equally possible. A Code byte 610 value of 138
25 corresponds to setting the Copy bit to 1, the Option class to 0, and the option number to
26 10, where 10 is one of the currently unused values for IP option numbers. The Option
27 number may alternatively be any currently unused IP option value.

28 Following the Code byte 610 in the encoding of the DAMP data packet IP options
29 field 520, the structure of a DAMP data packet header section 515 may follow a pattern
30 similar to that of the existing Strict Route IP Option, which is currently assigned to IP
31 Option number 137. See Comer, Douglas E., Internetworking with TCP/IP - Principles,
32 Protocols, and Architecture, 4th Ed., Vol. 1, Prentiss Hall, February 2000, pp.97-114.

1 The IP options field 520 is broken down into multiple 32-bit words, as shown in
2 Figure 6. The first word of the IP options field 520 contains the 8-bit (one byte) Code
3 byte 610, set to a value of 138 for DAMP data packets, and the Length byte 615, which
4 contains a value specifying the number of bytes in the IP options field 520 (including the
5 code and length bytes 610 and 615). This word is followed by multiple IP address fields
6 620, each containing four-byte IP addresses for each of the intended recipients of the
7 DAMP data packet 510, the remote network devices 240. There may be as many IP
8 address fields 620 as allowed for by the value specified in the Length byte 615. For
9 example, if the Length byte 615 indicates that the IP options field 520 includes 32 bytes,
10 then the IP options field 520 contains seven 32-bit (four-byte) IP address fields 620, in
11 addition to the first word containing the Code byte 610 and the Length byte 615. The
12 data packet header section 515 may also include two other fields of interest: a source IP
13 address field 522, specifying the IP address of the DAMP sending client 210, and a
14 destination IP address field 524, specifying the IP address of one of the remote network
15 devices 240. In an alternative embodiment, the IP address fields 620, the source IP
16 address field 522, and the destination IP address field 524 may each comprise multi-byte
17 IP addresses of a type other than four-byte addresses.

18 The data packet data section 525 is encoded following a standard User Datagram
19 Protocol ("UDP") IP data packet encoding scheme. As a connectionless protocol which,
20 like TCP, is layered on top of IP, UDP neither guarantees delivery nor does it require a
21 connection. As a result, it is lightweight and efficient, but all error processing and
22 retransmission must be taken care of by the application program. See Postel, Jon, User
23 Datagram Protocol, RFC 768, Network Information Center, SRI International, Menlo
24 Park, Calif., August 1980.

25 Figure 7 shows an exemplary path of one embodiment of a DAMP encoded data
26 packet 711 through the Internet to multiple destination remote network devices 240.
27 When a DAMP sending client 210 sends a DAMP data packet 711, it sets the source IP
28 address 522 contained in the data packet header section 525 to the IP address of the
29 DAMP sending client 210 itself, and the destination IP address 524 to any of the IP
30 addresses of the remote network devices 240 intended to receive the data packet. The
31 DAMP sending client 210 further encodes the data packet 711 with an IP Options field
32 520 containing the value for the DAMP IP Option, as described above, and which
33 contains a list of each of the desired destination remote network device 240 IP addresses.

1 The data packet data section 525, or UDP-encoded data portion of the data packet, thus
2 remains unchanged from that of a standard UDP data packet which may be sent to any
3 single recipient under existing IP transmission protocols.

4 When a router (or network switch) 715 receives a DAMP data packet 711, which
5 a router or switch may recognize from the DAMP IP Option encoded into the data packet
6 711, it will ignore the destination IP address 524, and instead examine the list of IP
7 addresses contained in the DAMP IP Options field 520. The router 715 then sends a copy
8 to each network interface 720 and 730 that contains at least one recipient as specified in
9 the address list in the IP Options field 520. However, before sending the data packets 721
10 and 731 on to a network interface, such as the network interfaces 720 and 730 shown,
11 each IP address in the IP Options field 520 list that is not intended to eventually receive
12 that copy of the data packet, i.e., is not found on a branch of that specific network
13 interface, is zeroed out. This prevents the creation of an infinite number of data packets
14 sent by two interconnected routers or switches addressing hosts existing across more than
15 one network interface. Furthermore, the switch or router 715 will set the destination IP
16 address 524 in the IP header section 515 to one of the non-zero entries remaining in the
17 recipient list (and the corresponding frame header destination hardware address). When a
18 DAMP data packet 721 or 731 is received with zeroed addresses in it, the zeroed
19 addresses are ignored. Entries are zeroed out rather than removing them in order that the
20 router or switch does not have to reformat the data packet.

21 Figure 8 is a diagram showing how a network of DAMP enabled routers and
22 switches would pass a DAMP data packet through several branches of the network.
23 Exemplary values for the relevant fields embedded in each successive copy of the DAMP
24 data packet are shown as it travels across each network segment. This figure shows how
25 in one embodiment of the invention, a network element 220, such as a router or IP switch
26 (815, 825, 835, 837), upon receiving a copy of the DAMP data packet 510 (811, 821,
27 831, 832) from the DAMP sending client 210 or from another network element (815, 825,
28 835, 837), then processes a copy of the received DAMP data packet (811, 821, 831, 832)
29 by zeroing those IP addresses embedded in the IP Options field 520 that are not directly
30 accessible below the network element (815, 825, 835, 837); setting the destination IP
31 address 524 in the IP header section 525 of the copy of the DAMP data packet (821, 822,
32 831, 832, 841, 842, 843) to the IP address of one of the non-zeroed remote destination
33 network devices (851, 852, 853, 854) embedded in the IP Options field 520; and routing

1 the modified copy of the DAMP data packet (821, 822, 831, 832, 841, 842, 843) to each
2 additional network element (825, 835, 837) or network device (851, 852, 853, 854) for
3 which there is a corresponding non-zeroed IP address listed in the embedded list of
4 multiple remote destination IP addresses.

5 When a remote network device 240 receives a DAMP data packet, it can process
6 it like any other UDP data packet. In fact, no DAMP software is necessary on the remote
7 network devices 240 beyond the normal UDP software necessary to receive the data
8 packet. In an alternate embodiment, the remote network device's 240 network element
9 220 could be programmed or otherwise adapted to be DAMP aware such that the network
10 element 220 could examine the list of addresses contained in the DAMP IP Options field
11 520 to determine if the DAMP data packet 510 was intended for that network element
12 220. This would remove the need for switches and routers to set the destination IP
13 address 524 to a recipient IP address (which may be impractical for network elements 220
14 to do without duplicating the data packet 510 and sending a copy of the data packet 510
15 for each recipient address to every network device 240 attached to it).

16 An alternative embodiment of DAMP data packet encoding may store destination
17 port numbers within the recipient list along with the IP addresses. Routers and switches
18 could then set the port numbers within the UDP header in the IP data section of the
19 DAMP data packet. This alternative embodiment may eliminate one potential limitation
20 of the DAMP encoding scheme whereby it is otherwise assumed that the destination port
21 number is identical for each recipient host.

22 The method 400, shown in Figure 4, is better understood in conjunction with the
23 diagram of the structure 500 of a DAMP data packet 510, as shown in Figure 5, and the
24 diagram of one embodiment of the structure of an IP options field 520 within a directly
25 addressed multicast protocol data packet 510, as shown in Figure 6.

26 The method 400 operates by embedding (step 420) multiple IP address fields 620
27 into the header section 515 of a DAMP data packet 510. Additionally, the method 400
28 provides a scheme (step 430) by which network elements 220, such as network routers or
29 switches, may access the list of multiple IP address fields 620. One embodiment of the
30 present invention places a uniquely formatted IP options field 520 in the IP header section
31 515 of the DAMP data packet 510. The IP options field 520, by way of a specific value
32 assigned to a Code byte 610 within the IP Options field 520, signifies that the data packet
33 510 is a DAMP data packet and thereby indicates that a list of multiple IP address fields

1 620 is embedded in the IP header section 515 of the DAMP data packet 510. The IP
2 options field 520 may then be read and translated (step 430) by network elements 220 and
3 the list of multiple IP address fields 620 embedded in the IP header section 515 of the
4 DAMP data packet 510 may then be accessed and read. The network elements 220 then
5 transmit (step 440) copies of the DAMP data packet 510 to each of the addresses in the
6 list of multiple IP address fields 620.

7 The present invention improves on existing multicast protocols in one respect by
8 not consuming bandwidth on a network segment to send the same data twice. This results
9 in lower network utilization and improved network performance. DAMP further
10 improves on existing multicast protocols by not requiring the overhead of setting up and
11 deleting multicast groups in the network switch, independent of the data packets, and by
12 not requiring a subscription service for the transmission of multicast data. This also
13 results in lower network utilization and improved performance.

14 For one embodiment of the present invention, a comparison may be made, outside
15 the unique DAMP addressing model and multicast behavior, between the network effects
16 of DAMP and the well known Unreliable Datagram Protocol ("UDP"), in contrast to the
17 Transmission Control Protocol ("TCP"). This analogy may be drawn where the
18 embodiment of the present invention does not require confirmation of receipt of the
19 DAMP data packets by the remote network devices to be sent back to the DAMP sending
20 client before continuing with transmission of additional DAMP data packets, as in UDP.
21 This embodiment substantially reduces the load on the network.

22 In another embodiment of the present invention, confirmation of receipt of the
23 DAMP data packets by the remote network devices is required to be sent back to the
24 DAMP sending client, similar to the manner in which TCP operates. This embodiment
25 requires bidirectional transmission of data, resulting in a correspondingly lesser decrease
26 in network traffic than the embodiment described above.

27 One potential use identified for DAMP is in networked storage systems. DAMP
28 allows data clients to send data directly to storage elements, and enables the data clients
29 to arrange the data in striped and mirrored configurations across numerous remote storage
30 elements, each with their own unique IP addresses, without the need for multicast
31 subscription management elements within the network infrastructure. DAMP allows the
32 client to efficiently and dynamically, with no unnecessary overhead, deliver data to the
33 appropriate remote storage elements.

1 The steps of the method 400 can be implemented with hardware or by execution
2 of programs, modules or scripts. The programs, modules or scripts can be stored or
3 embodied on one or more computer readable mediums in a variety of formats, such as
4 source code, object code or executable code, for example. The computer readable
5 mediums may include, for example, both storage devices and signals. Exemplary
6 computer readable storage devices include conventional computer system RAM (random
7 access memory), ROM (read only memory), EPROM (erasable, programmable ROM),
8 EEPROM (electrically erasable, programmable ROM), and magnetic or optical disks or
9 tapes. Exemplary computer readable signals, whether modulated using a carrier or not,
10 are signals that a computer system hosting or running the described methods can be
11 configured to access, including signals downloaded through the Internet or other
12 networks.

13 The terms and descriptions used herein are set forth by way of illustration only
14 and are not meant as limitations. Those skilled in the art will recognize that many
15 variations are possible within the spirit and scope of the invention as defined in the
16 following claims, and their equivalents, in which all terms are to be understood in their
17 broadest possible sense unless otherwise indicated